

## CLAIMS

What is claimed is:

- 1 1. A method for coordinating the writing of data items to persistent storage, the method  
2 comprising the steps of:  
3 maintaining within a first node a first queue for dirty data items that need to be  
4 written to persistent storage;  
5 maintaining within the first node a second queue for dirty data items that need to be  
6 written to persistent storage;  
7 moving entries from said first queue to said second queue when the dirty data items  
8 corresponding to the entries need to be transferred to a node other than said  
9 first node; and  
10 when selecting which data items to write to persistent storage, given priority to data  
11 items that correspond to entries in said second queue.
- 1 2. The method of Claim 1 wherein the step of moving entries includes moving an entry  
2 from said first queue to said second queue in response to a message received by said  
3 first node, wherein said message indicates that another node has requested the data  
4 item that corresponds to said entry.
- 1 3. A method for coordinating the writing of data items to persistent storage, the method  
2 comprising the steps of:  
3 maintaining a forced-write count for each of said data items;  
4 incrementing the forced-write count of a data item whenever the data item is written  
5 to persistent storage by one node for transfer of the data item to another node;  
6 selecting which dirty data items to write to persistent storage based on the write  
7 counts associated with the data items.

- 1 4. The method of Claim 3 further comprising the steps of:  
2 storing dirty data items that have forced-write counts above a certain threshold in a  
3 particular queue; and  
4 when selecting dirty data items to write to persistent storage, giving priority to data  
5 items stored in said particular queue.
- 1 5. A method of managing information about where to begin recovery after a failure, the  
2 method comprising the steps of:  
3 in a particular node of a multiple-node system, maintaining both  
4 a single-failure queue that indicates where within a recovery log to begin  
5 recovery after a failure of said node, and  
6 a multiple-failure queue that indicates where within said recovery log to begin  
7 recovery after a failure of said node and one or more other nodes in  
8 said multiple-node system;  
9 in response to a dirty data item being written to persistent storage, removing an entry  
10 for said data item from both said single-failure queue and said multiple-failure  
11 queue; and  
12 in response to a dirty data item being sent to another node of said multiple-node  
13 system without first being written to persistent storage, removing an entry for  
14 said data item from said single-failure queue without removing the entry for  
15 said data item from said multiple-failure queue.
- 1 6. The method of Claim 5 further comprising the step of sending the dirty data item to  
2 another node to allow removal of the entry from said single-failure queue without the  
3 other node requesting the dirty data item.
- 1 7. The method of Claim 5 further comprising the steps of:

2 after a single node failure, applying said recovery log beginning at a position in said  
3 recovery log associated with the single-failure queue; and  
4 after a multiple node failure, applying said recovery log beginning at a position in  
5 said recovery log associated with the multiple-failure queue.

1 8. The method of Claim 5 wherein:

2 said single-failure queue and said multiple-failure queue are implemented by a single  
3 combined queue; and  
4 the step of removing an entry for said data item from said single-failure queue  
5 without removing the entry for said data item from said multiple-failure queue  
6 includes marking an entry for said data item in said combined queue without  
7 removing the entry for said data item from said combined queue.

1 9. The method of Claim 5 wherein said single-failure queue and said multiple-failure  
2 queue are implemented as two separate queues.

1 10. A method for recovering after a failure, the method comprising the steps of:  
2 determining whether the failure involves only one node; and  
3 if the failure involves only said one node, then performing recovery by applying a  
4 recovery log of said node beginning at a first point in the recovery log; and  
5 if the failure involves one or more nodes in addition to said one node, then  
6 performing recovery by applying said recovery log of said node beginning at a  
7 second point in the recovery log;  
8 wherein said first point is different from said second point.

1 11. The method of Claim 10 wherein:

2 the first point is determined, at least in part, by which data items that were dirtied by  
3 said node reside in caches in other nodes; and

the second point is determined, at least in part, by which data items that were dirtied by said node have been persistently stored.

12. A method for recovering after a failure, the method comprising the steps of:  
if it is unclear whether a particular version of a data item has been written to disk,  
then performing the steps of  
without attempting to recover said data item, marking dirtied cached versions  
of said data item that would have been covered if said particular  
version was written to disk;  
when a request is made to write one of said dirtied cached versions to disk,  
determining which version of said data item is already on disk; and  
if said particular version of said data item is already on disk, then not writing  
said one of said dirtied cached versions to disk.

13. The method of Claim 12 further comprising the step of, if said particular version of said data item is not already on disk, then recovering said data item.

14. The method of Claim 12 further comprising the step of, if said particular version of said data item is already on disk, then informing nodes that contain said dirtied cached versions of the data item that said dirtied cached versions are covered by a write-to-disk operation.

15. A method for recovering a current version of a data item after a failure in a system that includes multiple caches, the method comprising the steps of:  
modifying the data item in a first node of said multiple caches to create a modified data item;

5 sending the modified data item from said first node to a second node of said multiple  
6 caches without durably storing the modified data item from said first node to  
7 persistent storage;  
8 after said modified data item has been sent from said first node to said second node  
9 and before said data item in said first node has been covered by a write-to-disk  
10 operation, discarding said data item in said first node; and  
11 after said failure, reconstructing the current version of said data item by applying  
12 changes to the data item on persistent storage based on merged redo logs  
13 associated with all of said multiple caches.

1 16. The method of Claim 15 further comprising the steps of:  
2 maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a  
3 locally-dirty checkpoint queue;  
4 wherein the globally-dirty data items associated with entries in the globally-dirty  
5 checkpoint queue are not retained until covered by write-to-disk operations;  
6 determining, for each cache, a checkpoint based on a lower of a first-dirtied time of  
7 the entry at the head of the locally-dirty checkpoint queue and the first-dirtied  
8 time of the entry at the head of the globally-dirty checkpoint queue; and  
9 after said failure, determining where to begin processing the redo log associated with  
10 each cache based on the checkpoint determined for said cache.

1 17. The method of Claim 15 further comprising the steps of:  
2 maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a  
3 locally-dirty checkpoint queue;  
4 wherein the globally-dirty data items associated with entries in the globally-dirty  
5 checkpoint queue are not retained until covered by write-to-disk operations;

maintaining, for each cache, a first checkpoint record for the locally-dirty checkpoint queue that indicates a first time, where all changes made to data items that are presently dirty in the cache prior to the first time have been recorded on a version of the data item that is on persistent storage;

maintaining, for each cache, a second checkpoint record for the globally-dirty checkpoint queue, wherein the second checkpoint record includes a list of data items that were once dirtied in the cache but have since been transferred out and not written to persistent storage; and

after said failure, determining where to begin processing the redo log associated with each cache based on the first checkpoint record and said second checkpoint record for said cache.

18. The method of Claim 17 wherein the step of processing the redo log comprises the steps of:
  - determining a starting position for scanning the redo log based on a lesser of
    - a position in the redo log as determined by the first checkpoint record and
    - the positions in the log as determined by the earliest change made to the list of the data items in the second checkpoint record; and
  - during recovery, for the portion of the redo log between the position indicated by the global checkpoint record to the position indicated by the local checkpoint record, considering for potential redo only those log records that correspond to the data items identified in the global checkpoint record.

19. A computer-readable medium carrying instructions for coordinating the writing of data items to persistent storage, the instructions comprising instructions for performing the steps of:

4 maintaining within a first node a first queue for dirty data items that need to be  
5 written to persistent storage;  
6 maintaining within the first node a second queue for dirty data items that need to be  
7 written to persistent storage;  
8 moving entries from said first queue to said second queue when the dirty data items  
9 corresponding to the entries need to be transferred to a node other than said  
10 first node; and  
11 when selecting which data items to write to persistent storage, given priority to data  
12 items that correspond to entries in said second queue.

1 20. The computer-readable medium of Claim 19 wherein the step of moving entries  
2 includes moving an entry from said first queue to said second queue in response to a  
3 message received by said first node, wherein said message indicates that another node  
4 has requested the data item that corresponds to said entry.

1 21. A computer-readable medium carrying instructions for coordinating the writing of  
2 data items to persistent storage, the instructions comprising instructions for  
3 performing the steps of:  
4 maintaining a forced-write count for each of said data items;  
5 incrementing the forced-write count of a data item whenever the data item is written  
6 to persistent storage by one node for transfer of the data item to another node;  
7 selecting which dirty data items to write to persistent storage based on the write  
8 counts associated with the data items.

1 22. The computer-readable medium of Claim 21 further comprising instructions for  
2 performing the steps of:  
3 storing dirty data items that have forced-write counts above a certain threshold in a  
4 particular queue; and

5 when selecting dirty data items to write to persistent storage, giving priority to data  
6 items stored in said particular queue.

1 23. A computer-readable medium carrying instructions for managing information about  
2 where to begin recovery after a failure, the instructions comprising instructions for  
3 performing the steps of:  
4 in a particular node of a multiple-node system, maintaining both  
5 a single-failure queue that indicates where within a recovery log to begin  
6 recovery after a failure of said node, and  
7 a multiple-failure queue that indicates where within said recovery log to begin  
8 recovery after a failure of said node and one or more other nodes in  
9 said multiple-node system;  
10 in response to a dirty data item being written to persistent storage, removing an entry  
11 for said data item from both said single-failure queue and said multiple-failure  
12 queue; and  
13 in response to a dirty data item being sent to another node of said multiple-node  
14 system without first being written to persistent storage, removing an entry for  
15 said data item from said single-failure queue without removing the entry for  
16 said data item from said multiple-failure queue.

1 24. The computer-readable medium of Claim 23 further comprising instructions for  
2 performing the step of sending the dirty data item to another node to allow removal of  
3 the entry from said single-failure queue without the other node requesting the dirty  
4 data item.

1 25. The computer-readable medium of Claim 23 further comprising instructions for  
2 performing the steps of:



3 after a single node failure, applying said recovery log beginning at a position in said  
4 recovery log associated with the single-failure queue; and  
5 after a multiple node failure, applying said recovery log beginning at a position in  
6 said recovery log associated with the multiple-failure queue.

1 26. The computer-readable medium of Claim 23 wherein:  
2 said single-failure queue and said multiple-failure queue are implemented by a single  
3 combined queue; and  
4 the step of removing an entry for said data item from said single-failure queue  
5 without removing the entry for said data item from said multiple-failure queue  
6 includes marking an entry for said data item in said combined queue without  
7 removing the entry for said data item from said combined queue.

1 27. The computer-readable medium of Claim 23 wherein said single-failure queue and  
2 said multiple-failure queue are implemented as two separate queues.

1 28. A computer-readable medium carrying instructions for recovering after a failure, the  
2 instructions comprising instructions for performing the steps of:  
3 determining whether the failure involves only one node; and  
4 if the failure involves only said one node, then performing recovery by applying a  
5 recovery log of said node beginning at a first point in the recovery log; and  
6 if the failure involves one or more nodes in addition to said one node, then  
7 performing recovery by applying said recovery log of said node beginning at a  
8 second point in the recovery log;  
9 wherein said first point is different from said second point.

1 29. The computer-readable medium of Claim 28 wherein:

2 the first point is determined, at least in part, by which data items that were dirtied by  
3 said node reside in caches in other nodes; and  
4 the second point is determined, at least in part, by which data items that were dirtied  
5 by said node have been persistently stored.

1 30. A computer-readable medium carrying instructions for recovering after a failure, the  
2 instructions comprising instructions for performing the steps of:  
3 if it is unclear whether a particular version of a data item has been written to disk,  
4 then performing the steps of  
5 without attempting to recover said data item, marking dirtied cached versions  
6 of said data item that would have been covered if said particular  
7 version was written to disk;  
8 when a request is made to write one of said dirtied cached versions to disk,  
9 determining which version of said data item is already on disk; and  
10 if said particular version of said data item is already on disk, then not writing  
11 said one of said dirtied cached versions to disk.

1 31. The computer-readable medium of Claim 30 further comprising instructions for  
2 performing the step of, if said particular version of said data item is not already on  
3 disk, then recovering said data item.

1 32. The computer-readable medium of Claim 30 further comprising instructions for  
2 performing the step of, if said particular version of said data item is already on disk,  
3 then informing nodes that contain said dirtied cached versions of the data item that  
4 said dirtied cached versions are covered by a write-to-disk operation.

1 33. A computer-readable medium carrying instructions for recovering a current version of  
2 a data item after a failure in a system that includes multiple caches, the instructions  
3 comprising instructions for performing the steps of:  
4 modifying the data item in a first node of said multiple caches to create a modified  
5 data item;  
6 sending the modified data item from said first node to a second node of said multiple  
7 caches without durably storing the modified data item from said first node to  
8 persistent storage;  
9 after said modified data item has been sent from said first node to said second node  
10 and before said data item in said first node has been covered by a write-to-disk  
11 operation, discarding said data item in said first node; and  
12 after said failure, reconstructing the current version of said data item by applying  
13 changes to the data item on persistent storage based on merged redo logs  
14 associated with all of said multiple caches.

1 34. The computer-readable medium of Claim 33 further comprising instructions for  
2 performing the steps of:  
3 maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a  
4 locally-dirty checkpoint queue;  
5 wherein the globally-dirty data items associated with entries in the globally-dirty  
6 checkpoint queue are not retained until covered by write-to-disk operations;  
7 determining, for each cache, a checkpoint based on a lower of a first-dirtied time of  
8 the entry at the head of the locally-dirty checkpoint queue and the first-dirtied  
9 time of the entry at the head of the globally-dirty checkpoint queue; and  
10 after said failure, determining where to begin processing the redo log associated with  
11 each cache based on the checkpoint determined for said cache.

1 35. The computer-readable medium of Claim 33 further comprising instructions for  
2 performing the steps of:  
3 maintaining, for each of said multiple caches, a globally-dirty checkpoint queue and a  
4 locally-dirty checkpoint queue;  
5 wherein the globally-dirty data items associated with entries in the globally-dirty  
6 checkpoint queue are not retained until covered by write-to-disk operations;  
7 maintaining, for each cache, a first checkpoint record for the locally-dirty checkpoint  
8 queue that indicates a first time, where all changes made to data items that are  
9 presently dirty in the cache prior to the first time have been recorded on a  
10 version of the data item that is on persistent storage;  
11 maintaining, for each cache, a second checkpoint record for the globally-dirty  
12 checkpoint queue, wherein the second checkpoint record includes a list of data  
13 items that were once dirtied in the cache but have since been transferred out  
14 and not written to persistent storage; and  
15 after said failure, determining where to begin processing the redo log associated with  
16 each cache based on the first checkpoint record and said second checkpoint  
17 record for said cache.

1 36. The computer-readable medium of Claim 35 wherein the step of processing the redo  
2 log comprises the steps of:  
3 determining a starting position for scanning the redo log based on a lesser of  
4 a position in the redo log as determined by the first checkpoint record and  
5 the positions in the log as determined by the earliest change made to the list of  
6 the data items in the second checkpoint record; and  
7 during recovery, for the portion of the redo log between the position indicated by the  
8 global checkpoint record to the position indicated by the local checkpoint

9 record, considering for potential redo only those log records that correspond to  
10 the data items identified in the global checkpoint record.